国际人工智能会议

Regret Bounds for Online Kernel Selection in Continuous Kernel Space

Xiao Zhang, Shizhong Liao, Jun Xu*, Ji-Rong Wen

Gaoling School of Artificial Intelligence, Renmin University of China College of Intelligence and Computing, Tianjin University

zhangx89@ruc.edu.cn



Introduction

2 Main Results



▲□▶★御★★書★★書★ 書 めへぐ

Vino	Thong	(DIIC)
AIdu	Znang	(NUC)

Offline / Online Learning



・ロマ・山田・山田・山中・日・シック

— Xi	20	7ŀ	an	σ (RΙ	(\mathbf{IC})
2 11	.uo .	~ .	uiii	ь \	1	υC,

Offline / Online Model Selection



Figure 1: Comparison between offline model selection and online model selection

Xiao Zhang (RUC)	AAAI 2021	December 19, 2020	4 / 22

Introduction

Necessity of Online Model Selection



Figure 2: Performances of Kernelized Online Gradient Decent (KOGD) and Kernel Perceptron (KP) using different Gaussian kernels for online classification

	•	다 에 관 에 관 에 관 에 관 에 관 .	9 Q (P
Xiao Zhang (RUC)	AAAI 2021	December 19, 2020	5 / 22

. —

Candidate Kernels

• Finite kernel set

containing a finite number of candidate kernels:

$$\widehat{\mathcal{H}}_N = \{\mathcal{H}_i\}_{i \in [N]},$$

where N is the number of candidate kernels.

• Continuous kernel space

containing continuously many candidate kernels:

$$\widehat{\mathcal{H}}_{\Omega} = \{ \mathcal{H}_{\sigma} \mid \sigma \in \Omega \},\$$

where Ω is the parameter interval of candidate kernels,

▲ □ ▶ ヵ 同 ヵ ヵ 日 ヵ ヵ 日

Regret of Online Kernel Selection

Definition 1 (Worse-Case Regret of Online Kernel Selection)

Given a continuous kernel space $\widehat{\mathcal{H}}_{\Omega}$ Assuming that $\{f_t\}_{t\in[T]} \subseteq \widehat{\mathcal{H}}_{\Omega}$ is a hypothesis sequence generated by an online kernel selection procedure after *T* rounds,

the worse-case regret (regret) of online kernel selection is defined as follows:

$$\mathcal{R}_{\text{reg}}\left(\{f_t\}_{t\in[T]},g\right) := \sum_{t=1}^{T} \left[\ell\left(f_t(\boldsymbol{x}_t),y_t\right) - \ell\left(g(\boldsymbol{x}_t),y_t\right)\right] = o(T),$$

where $g \in \widehat{\mathcal{H}}_{\Omega}$ is the competing hypothesis, which is typically defined as the best hypothesis in hindsight:

$$g = f^* := \operatorname*{arg\,min}_{f \in \widehat{\mathcal{H}}_{\Omega}} \sum_{t=1}^T \ell\left(f(\boldsymbol{x}_t), y_t\right).$$

▲日▼★聞★★回★★回★ 団・⊘へ⊙

Introduction

New Challenges of Online Kernel Selection

- (1) High computational complexities
 - Quadratic time complexity [Singh and Príncipe, 2011]
 - Linear space complexity [Chen et al., 2016]
- (2) Lack of theoretical guarantees
 - Regret bounds dependent on the number of kernels [Yang et al., 2012]
 - Expected regret bounds in continuous kernel space [Zhang and Liao, 2020]

イロト (日本) イヨト

DQ P

Main Contribution

Table 1: Comparison with the existing online kernel selection approaches

	Computa	tional complexiti	Theoretical guarantees		
	Time (per round)	Time (overall)	Space	#Candidate	Regret bound
Existing	Linear	Quadratic	Linear	Finite	Square root
Proposed	Constant	Linear	Constant	Infinite	Square root

Learning Framework



Figure 3: Online kernel selection using time-varying hypothesis sketching

Xiao Zhang (RUC)

900

◆□▶→□→→□=→→□=→→□=

Time-Varying Hypothesis Sketch

• Time-varying hypothesis sketch at round *t*

$$f_{\sigma_t,t}(\cdot) = \langle \boldsymbol{\omega}^{(t)}, \boldsymbol{\psi}^{(t)}_{\sigma_t}(\cdot) \rangle, \quad \sigma_t \in \Omega,$$

where $\mathcal{V}_t = {\{\tilde{\boldsymbol{x}}_i\}}_{i=1}^{|\mathcal{V}_t|} \subseteq \mathcal{X}$ is a buffer.

• Basis vector and its corresponding weight vector

$$\boldsymbol{\psi}_{\sigma_{t}}^{(t)}(\cdot) = \begin{bmatrix} \kappa_{\sigma_{t}}(\cdot, \tilde{\boldsymbol{x}}_{1}), \dots, \kappa_{\sigma_{t}}(\cdot, \tilde{\boldsymbol{x}}_{|\mathcal{V}_{t}|}) \end{bmatrix}^{\mathsf{T}}, \quad \tilde{\boldsymbol{x}}_{i} \in \mathcal{V}_{t}, \\ \boldsymbol{\omega}^{(t)} = \begin{bmatrix} \omega_{1}^{(t)}, \omega_{2}^{(t)}, \dots, \omega_{|\mathcal{V}_{t}|}^{(t)} \end{bmatrix}^{\mathsf{T}} \in \mathbb{R}^{|\mathcal{V}_{t}|}.$$

Ŧ

メロト 日間 ちょうちょう

Learning Framework



Figure 4: Online kernel selection using time-varying hypothesis sketching

< □ ▶

900

12/22

Xiao Zhang (RUC)AAAI 2021December 19, 2020

Xiao Zhang (RUC)

Two Online Kernel Selection Categories



Figure 5: Two Online Kernel Selection Categories at round *t*: OKS-SPT (left) denotes the online kernel selection by **selection-post-training**; OKS-TPS (right) denotes the online kernel selection by **training-post-selection**.

AAAI 2021

				Ľ	ece	eml	ber	19	, 20	20	13 /	22
•	¢		• 🗗		Ē	5.		Ē	8		Э с	$\langle \mathcal{O} \rangle$

Two Online Kernel Selection Categories

Category 1: Online Kernel Selection by Selection-Post-Training (OKS-SPT)

Require: the continuous kernel space \mathcal{K}_{Ω} , initial kernel κ_{σ_1} 1: Initialize the weight vector $\boldsymbol{\omega}^{(1)} = \mathbf{0}$ 2: for t = 1, ..., T do Compute the hypothesis sketch $f_{\sigma_t,t}(\cdot) = \langle \boldsymbol{\omega}^{(t)}, \boldsymbol{\psi}^{(t)}_{\sigma_t}(\cdot) \rangle$ 3: Predict $\hat{y}_t = \operatorname{sgn}(f_{\sigma_t,t}(\boldsymbol{x}_t))$ for classification or 4: $\hat{y}_t = f_{\sigma_t,t}(\boldsymbol{x}_t)$ for regression Maintain the buffer $\mathcal{V}_{t+1} = \text{BUFFERMAINTENANCE}(\mathcal{V}_t, \boldsymbol{z}_t)$ and obtain $f_{\sigma_t,t}^{\mathrm{ma}} = \langle \boldsymbol{\omega}^{(t)}, \boldsymbol{\psi}_{\sigma_t}^{(t+1)}(\cdot) \rangle$ Update the weight vector 6: $\boldsymbol{\omega}^{(t+1)} = \text{WEIGHTUPDATING}(f_{\sigma_t,t}^{\text{ma}}, \boldsymbol{z}_t)$ and obtain $f_{\sigma,t+1}^{\mathrm{ma}} = \langle \omega^{(t+1)}, \psi^{(t+1)}_{\sigma,t}(\cdot) \rangle$ if the buffer changes then 7: 8: Select $\kappa_{\sigma_{t+1}} = \text{KERNELSELECTION}(\kappa_{\sigma_t}, \boldsymbol{z}_t)$ from \mathcal{K}_{Ω} and obtain $f_{\sigma_{t+1},t+1}$ 9: else 10: $\kappa_{\sigma_{t+1}} = \kappa_{\sigma_t}$ 11: end if 12: end for

• • • • • • • • • • • • •

DQ P

Two Online Kernel Selection Categories

Category 2: Online Kernel Selection by Training-Post-Selection (OKS-TPS)

Require: the continuous kernel space \mathcal{K}_{Ω} , initial kernel κ_{σ_1} 1: Initialize the weight vector $\boldsymbol{\omega}^{(1)} = \mathbf{0}$ 2: for t = 1, ..., T do Compute the hypothesis sketch $f_{\sigma_t,t}(\cdot) = \langle \boldsymbol{\omega}^{(t)}, \boldsymbol{\psi}^{(t)}_{\sigma_t}(\cdot) \rangle$ 3: Predict $\hat{y}_t = \operatorname{sgn}(f_{\sigma_t,t}(\boldsymbol{x}_t))$ for classification or 4: $\hat{y}_t = f_{\sigma_t,t}(\boldsymbol{x}_t)$ for regression Select a kernel $\kappa_{\sigma_{t+1}} = \text{KERNELSELECTION}(\kappa_{\sigma_t}, \boldsymbol{z}_t)$ 5: from \mathcal{K}_{Ω} Maintain the buffer 6: $\mathcal{V}_{t+1} = \text{BUFFERMAINTENANCE}(\mathcal{V}_t, \boldsymbol{z}_t)$ and obtain $f_{\sigma_{t+1},t}^{\mathrm{ma}} = \langle \boldsymbol{\omega}^{(t)}, \boldsymbol{\psi}_{\sigma_{t+1}}^{(t+1)}(\cdot) \rangle$ Update the weight vector 7: $\boldsymbol{\omega}^{(t+1)} = \text{WEIGHTUPDATING}(f_{\sigma_{t+1},t}^{\text{ma}}, \boldsymbol{z}_t) \text{ and obtain}$ $f_{\sigma_{t+1},t+1}$ 8: end for

JOC.

◆□▶ ○冊 ○ ○臣 ○ ○臣 ○

Main Results

Our Theory: Regret Analysis

Theorem 2 (Regret of OKS-SPT)

Let \mathcal{K}_{Ω} be a continuous kernel space that contains Gaussian kernel functions, $\{f_{\sigma_t,t}\}_{t=1}^T \subseteq \widehat{\mathcal{H}}$ be the hypothesis sketch sequence generated by OKS-SPT. Define $\widehat{D}(f,g) = \max_{t \in [T]} |f(\mathbf{x}_t) - g(\mathbf{x}_t)|$, where $f, g \in \widehat{\mathcal{H}}$. Assume $C_{\max} = \max_{i,j \in [T]} ||\mathbf{x}_i - \mathbf{x}_j||^2$ and $R = \sup_{f \in \widehat{\mathcal{H}}} ||f||_{\widehat{\mathcal{H}}}$, then there exists a constant C > 0 such that

$$\widehat{\operatorname{Reg}}_{T}(\{f_{\sigma_{t},t}\}_{t=1}^{T},\overline{f}^{*}) \leq \underbrace{U_{1}}_{U_{1}} + \underbrace{U_{2}}_{U_{2}} + \underbrace{U_{3}}_{U_{3}}$$

Optimization Estimation Approximation

where $U_2 = \widehat{D}(\overline{f}^*_{\sigma_t}, f^*_{\sigma_t})$,

$$U_{1} = 2RC\sqrt{\mu} \left[T - O(\ln T)\right] + \frac{R^{2}}{2\eta_{f}} + \frac{\left(C\sqrt{\mu} + 1\right)^{2}}{2}\eta_{f}T,$$

$$U_{3} = 2C_{\max}\frac{\sigma_{\max}}{\sigma_{\min}^{3}}\widehat{D}(\bar{f}_{\sigma_{t}}^{*}, f_{\sigma_{t}, t+1}) + \frac{(\bar{\sigma}^{*})^{2}}{2\eta_{\sigma}} + \frac{\left[\sigma_{\min}^{-3}C_{\max}\widehat{D}(\bar{f}_{\sigma_{t}}^{*}, f_{\sigma_{t}, t+1}) + L\right]^{2}}{2}\eta_{\sigma}T.$$

Our Theory: Regret Analysis

- (1) Selection-post-training (OKS-SPT): $O(\sqrt{T})$ 后悔界
 - Stepsize: $\eta_f, \eta_\sigma = O(1/\sqrt{T})$
 - Threshold value for buffer: $\nu = O(1/T)$
 - Hypothesis value: $O(1/\sqrt{T})$
- (2) Selection-post-training (OKS-TPS): $O(\sqrt{T})$ 后悔界
 - Stepsize: $\eta_f, \eta_\sigma = O(1/\sqrt{T})$
 - Threshold value for buffer: $\nu = O(1/T)$
 - Conditions for kernel matrix:

$$\det \mathbf{K}_{\bar{\sigma}^*,\mathcal{V}_t \cup \{\mathbf{x}_t\}} / \det \mathbf{K}_{\bar{\sigma}^*,\mathcal{V}_t} = O(\mu),$$

$$\overset{||}{(\mathbf{K}_{\bar{\sigma}^*,\mathcal{V}_t})^{\dagger} \psi_{\bar{\sigma}^*}^{(t)}(\mathbf{x}_t) - (\mathbf{K}_{\sigma_{t+1},\mathcal{V}_t})^{\dagger} \psi_{\sigma_{t+1}}^{(t)}(\mathbf{x}_t)^{||} = O(\mu).$$

▲ □ ▶ ヵ 同 ゃ ヵ 三 ゃ ぉ

 $\mathcal{A} \mathcal{A} \mathcal{A}$

Our Theory: Comparison of Theoretical Results

Approach	Com	putational Comp	Regret Guarantees		
	#Updates	Time (overall)	Space	Candidate	Regret bound
OKS	Т	$O(T^2 + NT)$	O(T)	Finite	$O\left(\sqrt{N(\ln N)T}\right)$
MS-FTPL	Т	$O\left(NT^2\right)$	O(NT)	Finite	$O\left(\sqrt{NT\ln T}\right)^{2}$
OKL-GD	Т	$O(T^2)$	O(T)	Continuous	_
RRF	Т	$O(T^2)$	O(T)	Continuous	_
OKS-SPT	$\ln(T)$	$O((\ln T)^2 T)$	$O((\ln T)^2)$	Continuous	$O\left(\sqrt{T}\right)$
OKS-TPS	Т	$O((\ln T)^2 T)$	$O((\ln T)^2)$	Continuous	$O\left(\sqrt{T}\right)$

Table 2: Comparison among online kernel selection approaches

• *T*: the number of rounds.

• *N*: the number of candidate kernels.

▲□▶ 4冊 4 4 글 4 4 글

Ŧ

Empirical Verification

Table 3: Performances of OKL-GD, OKS, RRF and the proposed OKS-SPT, OKS-TPS for online classification w.r.t. the mistake rate (%) and the running time (s).

Algorithm	germa	n	spambas	se	mushroo	mushrooms		
Algorium	Mistake rate	Time	Mistake rate	Time	Mistake rate	Time		
OKL-GD	34.960 ± 1.518	0.194	36.031 ± 0.421	6.700	4.226 ± 0.910	37.230		
OKS	42.320 ± 1.307	0.226	34.355 ± 0.372	4.083	9.441 ± 0.282	9.787		
RRF	31.140 ± 0.114	0.372	44.961 ± 0.820	4.767	16.166 ± 0.964	21.750		
OKS-SPT	29.920 ± 0.286	0.244	$\textbf{28.436} \pm \textbf{0.213}$	2.533	6.585 ± 0.246	4.240		
OKS-TPS	$\textbf{29.760} \pm \textbf{0.270}$	0.296	28.450 ± 0.188	2.590	$\textbf{3.139} \pm \textbf{0.481}$	6.910		
Algorithm	a9a		w7a		ijcnn1	ijcnn1		
Algorium	Mistake rate	Time	Mistake rate	Time	Mistake rate	Time		
OKL-GD	23.936 ± 0.008	321.525	2.975 ± 0.062	857.268	9.575 ± 0.012	134.880		
OKS	23.617 ± 0.127	1053.420	7.637 ± 0.024	943.855	9.578 ± 0.184	618.520		
RRF	23.931 ± 0.001	152.265	2.978 ± 0.004	674.735	9.574 ± 0.001	39.290		
OKS-SPT	$\textbf{20.368} \pm \textbf{0.659}$	39.360	2.675 ± 0.023	94.530	9.478 ± 0.003	33.860		
OKS-TPS	22.379 ± 0.192	48.815	$\textbf{2.631} \pm \textbf{0.010}$	96.395	$\textbf{9.440} \pm \textbf{0.002}$	35.165		

▲□▶★聞★★臣★★臣★ 臣 めんの

Conclusion

- Two online kernel selection categories via time-varying hypothesis sketching.
- Meet the new challenges of online kernel selection.
- A time-varying sketching approach to online model selection.

1

▲□▶ → □ → → 三 → →

Main References

[Chen et al., 2016] Chen, B., Liang, J., Zheng, N., and Príncipe, J. C. (2016). Kernel least mean square with adaptive kernel size. *Neurocomputing*, 191:95–106.

[Singh and Príncipe, 2011] Singh, A. and Príncipe, J. C. (2011). Information theoretic learning with adaptive kernels. *IEEE Transactions on Signal Processing*, 91(2):203–213.

[Yang et al., 2012] Yang, T., Mahdavi, M., Jin, R., Yi, J., and Hoi, S. C. (2012).
 Online kernel selection: Algorithms and evaluations.
 In *Proceedings of the 26th AAAI Conference on Artificial Intelligence*, pages 22–26.

[Zhang and Liao, 2020] Zhang, X. and Liao, S. (2020).
 Hypothesis sketching for online kernel selection in continuous kernel space.
 In Proceedings of the 29th International Joint Conference on Artificial Intelligence, pages 2498–2504.

国际人工智能会议 AAAI 2021论文北京预讲会

THANKS

2020.12.19